# Progressive Join Algorithms Considering User Preference

Mengsu Ding, **Shimin Chen**\*

Institute of Computing Technology
Chinese Academy of Sciences
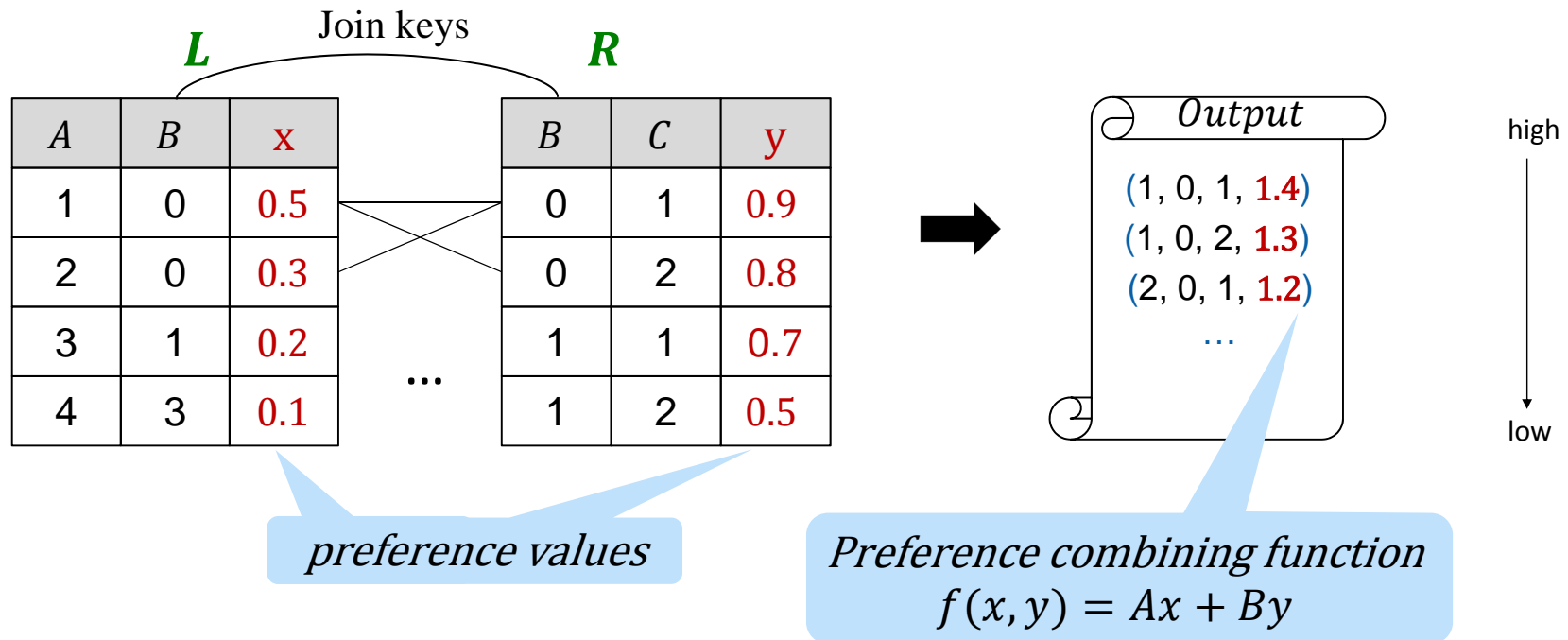Beijing, China

Nantia Makrynioti, Stefan Manegold

CWI (Centrum Wiskunde & Informatica)
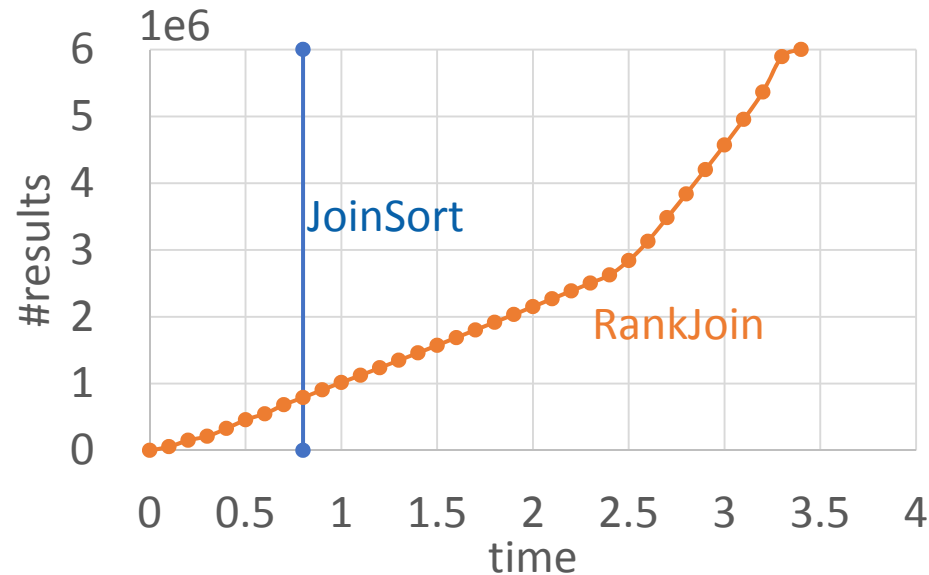Amsterdam, the Netherlands

# Preference-aware Progressive Join

- Progressive query processing for exploratory data analysis

- Return join results ordered according to preference



Join keys

*L*        *R*

| *A* | *B* | x |
|---|---|---|
| 1 | 0 | 0.5 |
| 2 | 0 | 0.3 |
| 3 | 1 | 0.2 |
| 4 | 3 | 0.1 |

| *B* | *C* | y |
|---|---|---|
| 0 | 1 | 0.9 |
| 0 | 2 | 0.8 |
| 1 | 1 | 0.7 |
| 1 | 2 | 0.5 |

*Output*

(1, 0, 1, **1.4**)
(1, 0, 2, **1.3**)
(2, 0, 1, **1.2**)
...

high

low

*preference values*

*Preference combining function*
$f(x, y) = Ax + By$

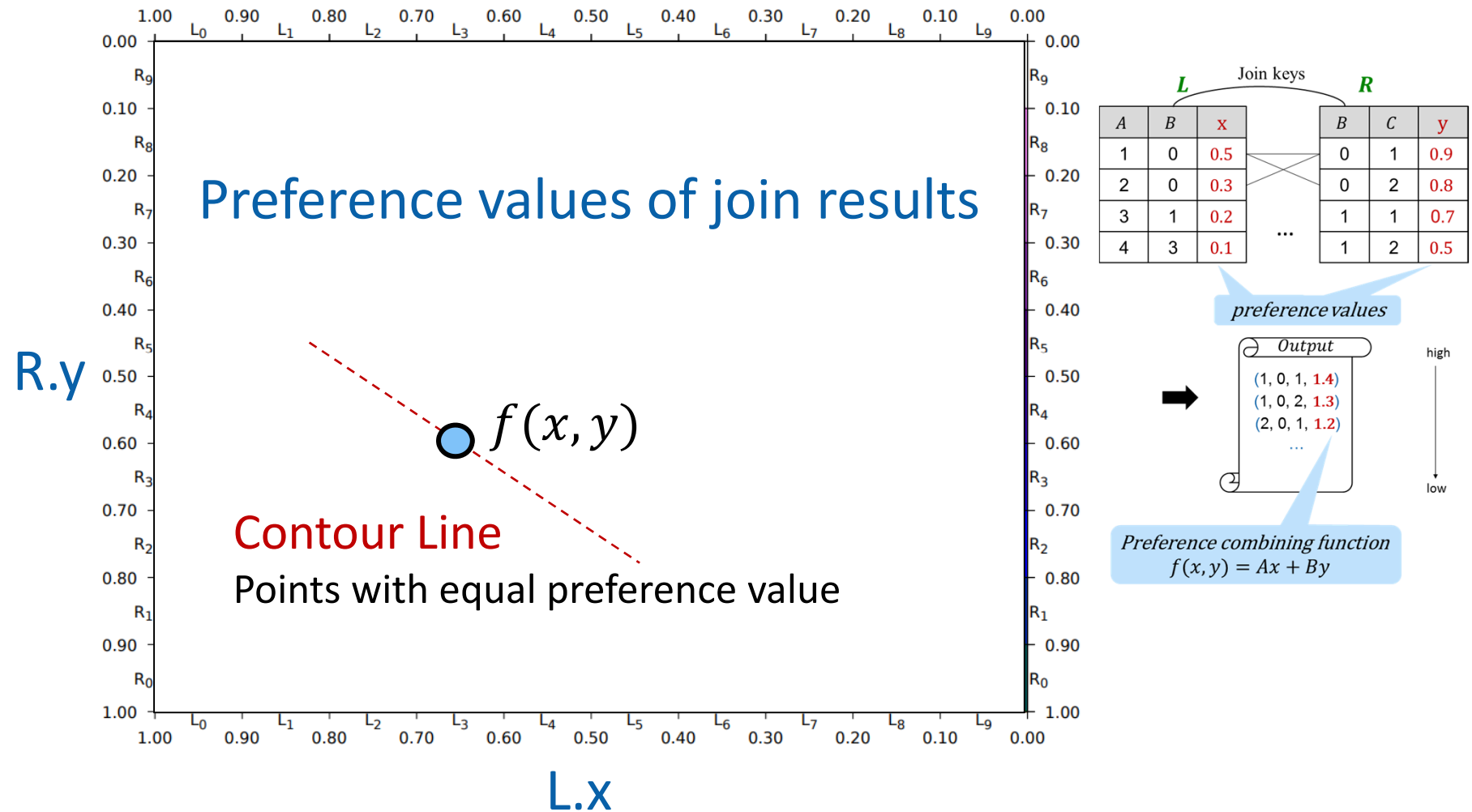- Goal: Fast early results & Fast full results
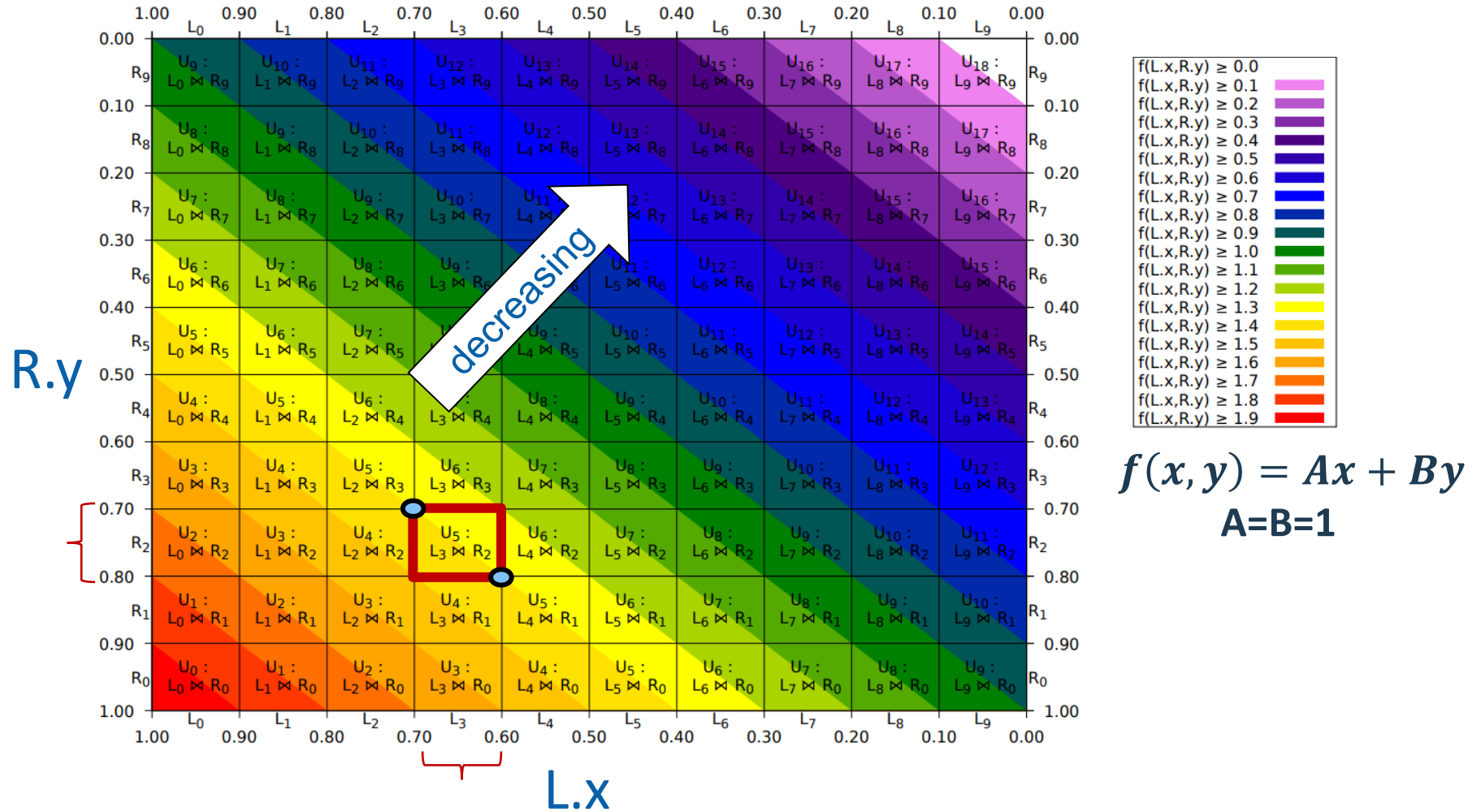
# Problems of Existing Solutions



- Blocking approach: Join + Sort
  - ❑ **No early results**

- Extending top-$k$ join algorithms:
  - ❑ e.g., RankJoin [VLDB'03] symmetric hash join + priority queue
  - ❑ **Slow full results**
  - ❑ **Significant sorting overhead** ⟵ We want to reduce or eliminate sorting overhead
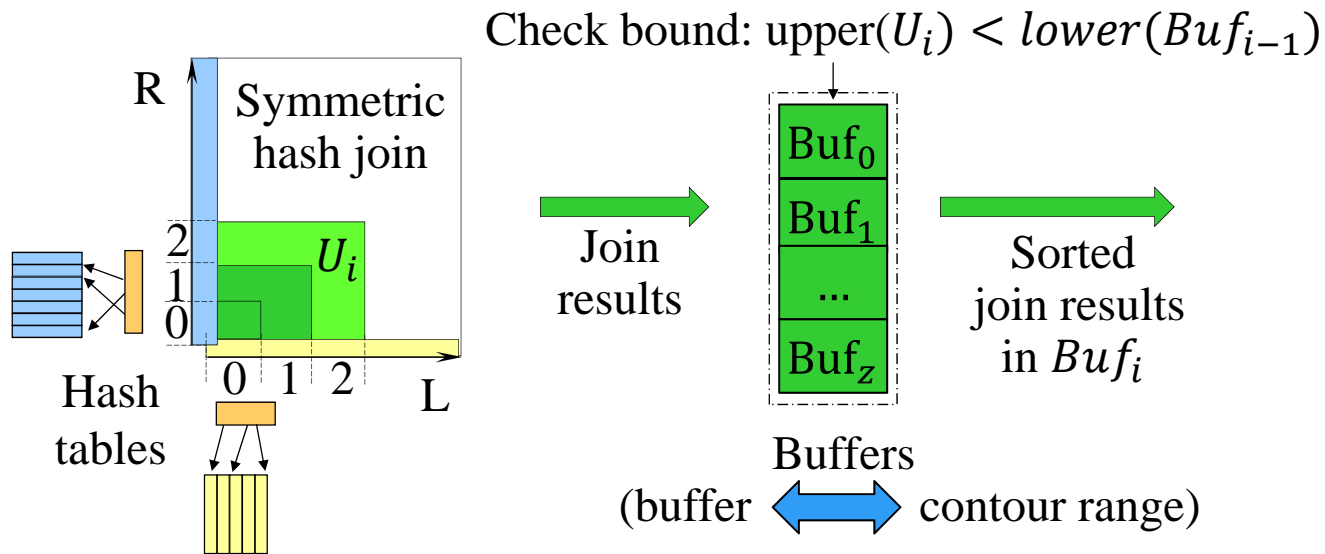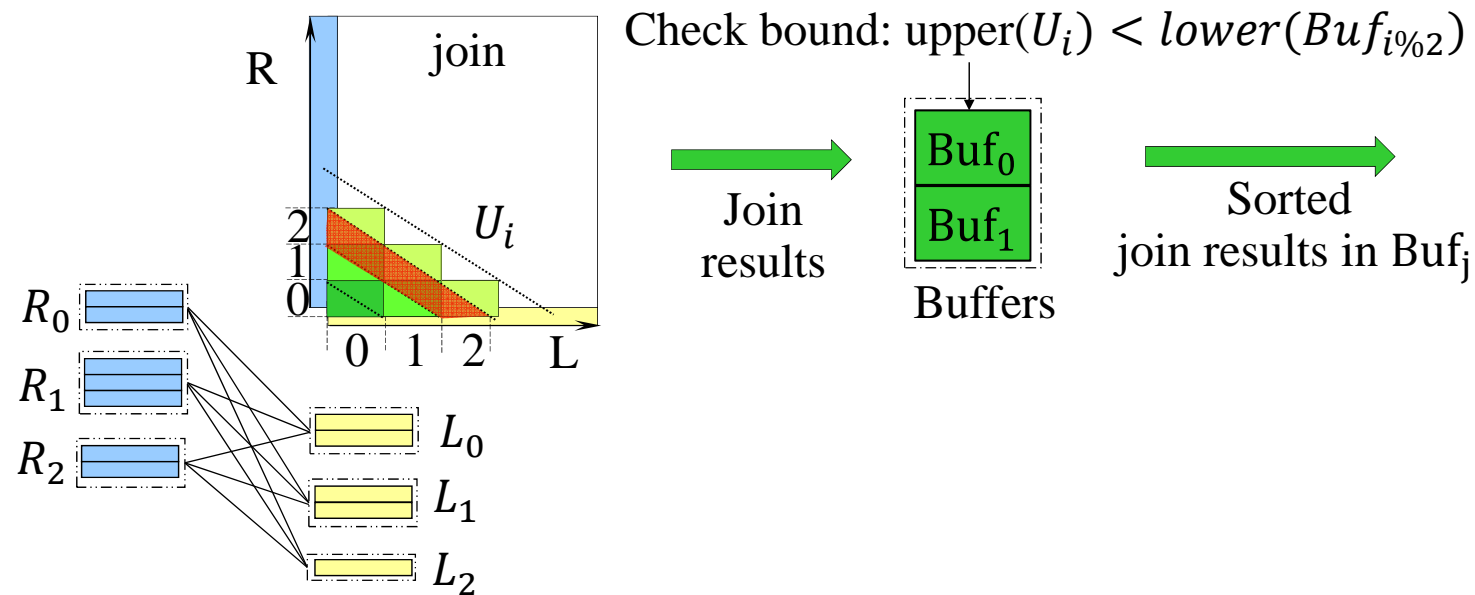
# Our Idea: Exploiting Contour Lines



Preference values of join results

R.y

$f(x, y)$

Contour Line
Points with equal preference value

L.x

# Our Idea: Exploiting Contour Lines



$$f(x, y) = Ax + By$$

A=B=1

# Inputs Follow Contour Lines (Algorithm1: CJp**I**)

Check bound: $upper(U_i) < lower(Buf_{i-1})$

R

Symmetric hash join

$U_i$

Hash tables

Join results

$Buf_0$
$Buf_1$
...
$Buf_z$

Sorted join results in $Buf_i$

Buffers

(buffer ⟷ contour range)

- Avoid sorting across buffers
- Only need to sort within each buffer

# Both Inputs and Outputs Follow Contour Lines (Algorithm2: CJp**B**)



Check bound: $upper(U_i) < lower(Buf_{i\%2})$

- Reduce intermediate result size
- But join cost may be higher

# Relaxation to Remove Sorting

- Relaxation of Order: **no intra-buffer sorting**

  ❑ Relaxation: tuples $t_i$ and $t_j$ are regarded as in order if
  $$\left| t_i.pval - t_j.pval \right| \leq \epsilon$$
  ❑ Judiciously select input intervals

### 4 Variants of Contour Joins

| Variants | Follow Contour Lines | | Relaxation |
| :---: | :---: | :---: | :---: |
| | Join Inputs | Join Results | |
| CJpI | ✓ | | |
| CJpB | ✓ | ✓ | |
| CJrI | ✓ | | ✓ |
| CJrB | ✓ | ✓ | ✓ |

p: precise, r: relaxed;    I: Inputs, B: Both inputs & outputs

# Experiments

- Data Set: Based on TPC-H Lineitem and Partsupp

  ❑ Preference values based on: l_discount, ps_availqty

  ❑ 3 datasets (can fit into main memory)

| Scale Factor | Datasets | #Inputs | #Outputs |
|---|---|---|---|
| $m_1$ | $datasets\#1$ | **fixed** | **increasing** |
| $m_2$ | $datasets\#2$ | **increasing** | **fixed** |
| $m_3$ | $datasets\#3$ | **increasing** | |

Join key

**PARTSUPP (PS_)**
SF*800,000

| PARTKEY |
| SUPPKEY |
| AVAILQTY |
| SUPPLYCOST |
| COMMENT |

**LINEITEM (L_)**
SF*6,000,000

| ORDERKEY |
| PARTKEY |
| SUPPKEY |
| LINENUMBER |
| QUANTITY |
| EXTENDEDPRICE |
| DISCOUNT |
| TAX |

*preference values*

source: TPCH spec

- Query

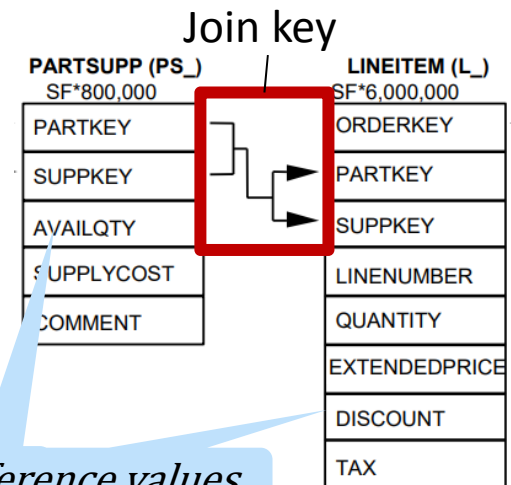$$f(x, y) = Ax + By$$

select   L.key, $f(L.pval, PS.pval)$ as score
from     Lineitem as L, Partsupp as PS
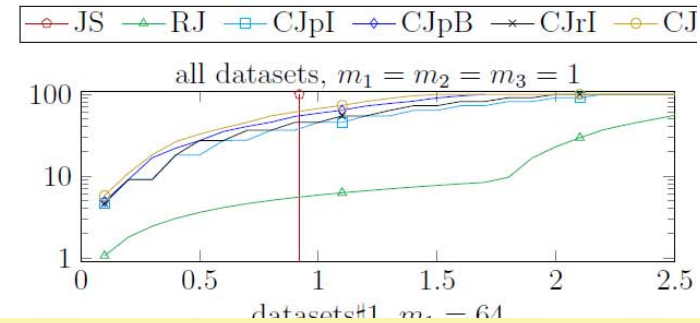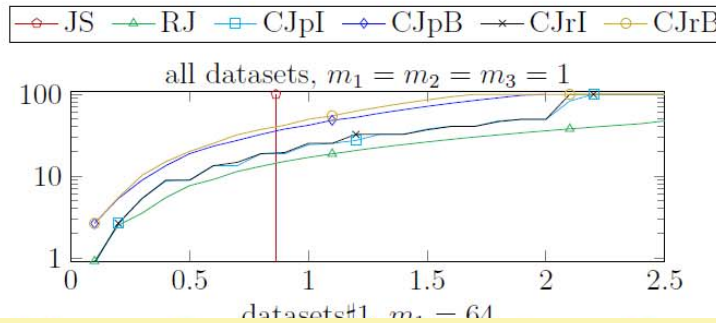where   L.key = PS.key
order by score ASC
progressively

- Machine

  ❑ Intel Core i7-4770 CPU @3.40GHz (8MB cache) and 32GB memory

  ❑ 64-bit Ubuntu 16.04 LTS with 4.15.0-62-generic Linux kernel

# Overall Results



all datasets, $m_1 = m_2 = m_3 = 1$

Legend: JS, RJ, CJpI, CJpB, CJrI, CJrB

- **Compared to RJ (RankJoin)**
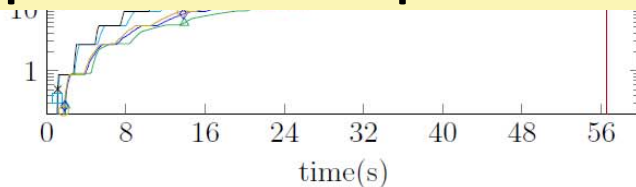  - ❑ 1% early results
    - best **precise** contour join: up to 7x improvements
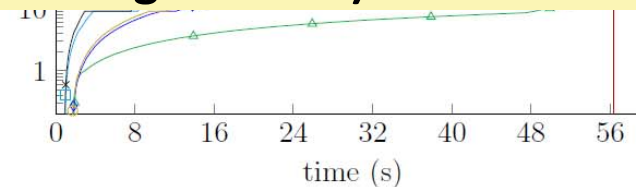    - best **relaxed** contour join: up to 14x improvements
  - ❑ Full results
    - best **precise** contour join: up to 10.6x improvements
    - best **relaxed** contour join: up to 39.4x improvements
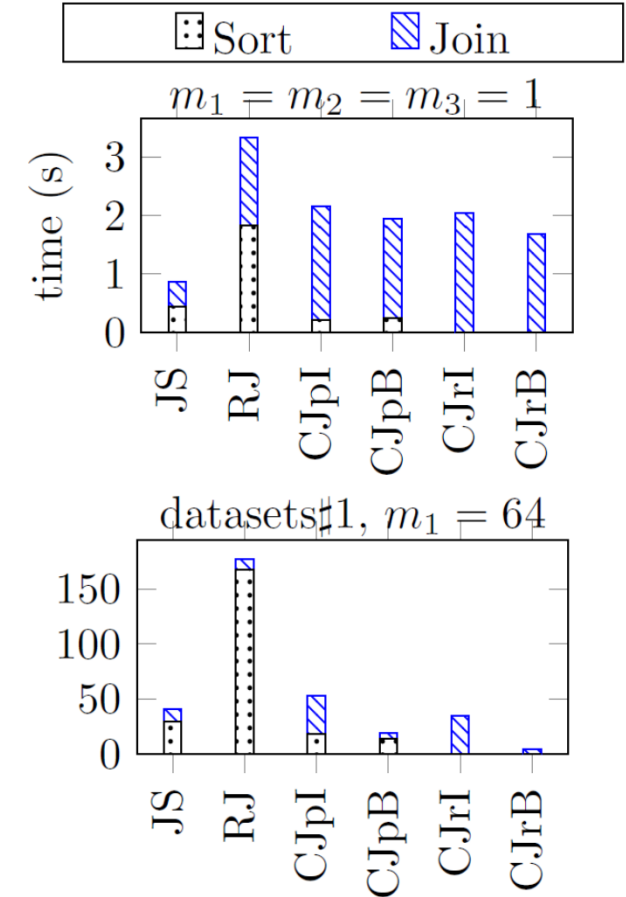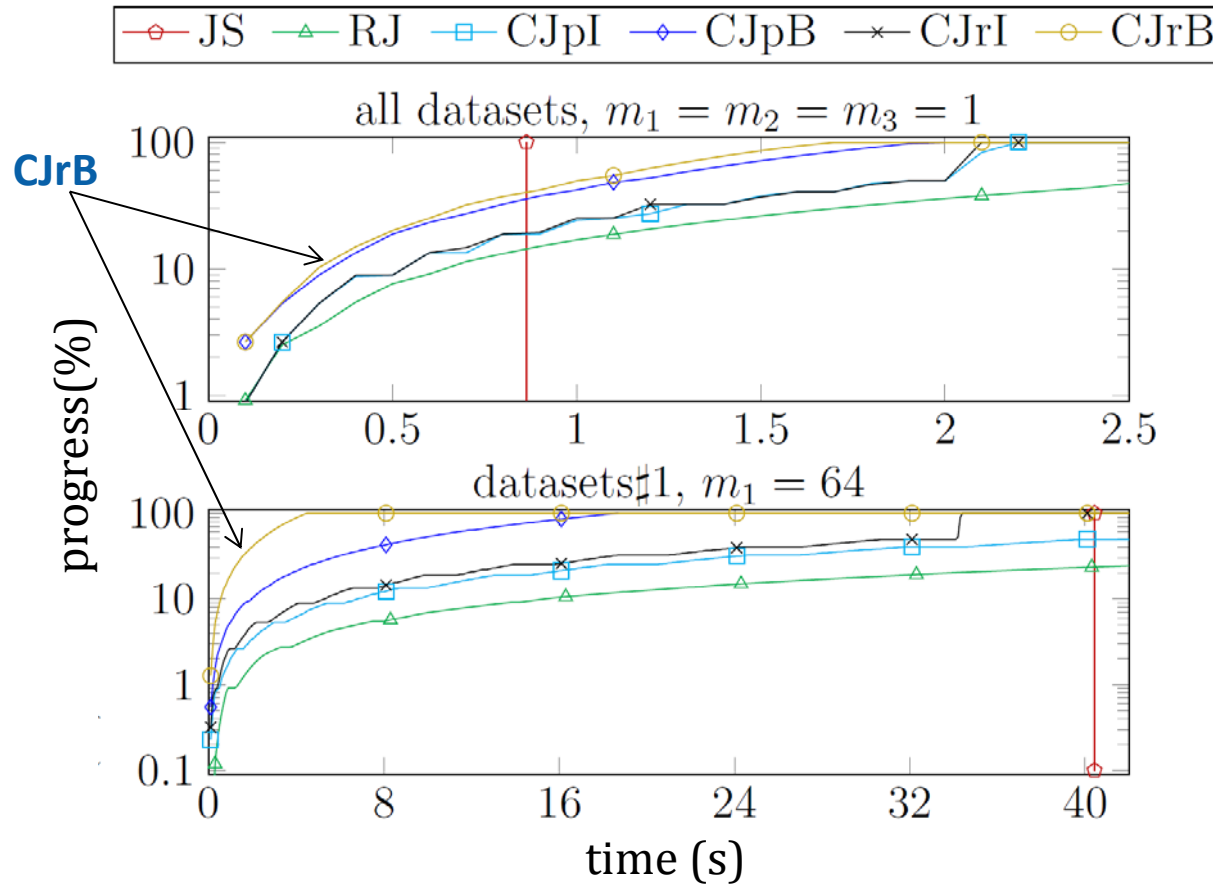- **Comparable or better perfomance to JS (blocking JoinSort)**

(a) $f(x, y) = x + y$, $p_L = 200$, $p_R = 200$

(b) $f(x, y) = 10x + y$, $p_L = 2000$, $p_R = 200$

# Fix Input Size, Increase Join Result Size



- As join result size increases, the fraction of sorting increases
- RankJoin becomes very poor
- CJrB is the best performing algorithm

# Conclusion

- Preference-aware joins in progressive query processing

- Idea: exploit contour lines in the join result space

- ContourJoin: a promising solution
  - ❑ Faster early and full results generation (vs. RankJoin)
  - ❑ Good total result generation performance (vs. JoinSort)

  - ❑ More in the paper
    - – Algorithms and proofs
    - – Extensive experimental results
    - – Discussion on preference combining functions, unsorted inputs, multi-way joins